# Statistical Problems in a Paper on Variation In Cancer Risk Among Tissues, and New Discoveries

LEE ALTENBERG

The KLI Institute, Lee.Altenberg@kli.ac.at.

January 20, 2015

## Abstract

Tomasetti and Vogelstein (2015) collected data on 31 different tissue types and found a correlation of $0.8$ between the logarithms of the incidence of cancer (LCI), and the estimated number of stem cell divisions in those tissues (LSCD). Some of their conclusions however are statistically erroneous. Their excess risk score, "ERS" ($\log_{10}$ LCI $\times \log_{10}$ LSCD), is non-monotonic under a change of time units for the rates, which renders meaningless the results derived from it, including a cluster of 22 "R-tumor" types for which they conclude, "primary prevention measures are not likely to be very effective". Further, $r = 0.8$ is consistent with the three orders of magnitude variation in other unmeasured factors, leaving room for the possibility of primary prevention if such factors can be intervened upon. Further exploration of the data reveals additional findings: (1) that LCI grows at approximately the square root of LSCD, which may provide a clue to the biology; (2) among different possible combinations of the primary data, the one maximizing the correlations with LCI is almost precisely the formula used by Tomasetti and Vogelstein to estimate LSCD, giving support to stem cell divisions as an independent factor in carcinogenesis, while not excluding other such factors.

*Key words: cancer incidence, stem cells, correlation, tumor types, linear regression, somatic mutation, fallacy, symbolic regression, prevention*

## Introduction

A paper recently published in *Science*, "Variation in cancer risk among tissues can be explained by the number of stem cell divisions" by Tomasetti and Vogelstein (2015), investigates the relationship between the number of stem cell divisions in various tissue types and the incidence of tumors in them. They find a strong correlation, $0.8$, between the logarithms of these two values (both Spearman's rank correlation $\rho$ and Pearson's linear correlation $r$). They use the product of these logarithms, which they call the "extra risk score" (ERS) to classify tumor types

into "R-tumors", where "stochastic factors, presumably related to errors during DNA replication, most strongly appear to affect their risk," and "D-tumors", where "deterministic factors such as environmental mutagens or hereditary predispositions strongly affect their risk". On this basis, they give a prognosis for the likely success of interventions to prevent these tumor types:

> These results could have important public health implications. One of the most promising avenues for reducing cancer deaths is through prevention. How successful can such approaches be? The maximum fraction of tumors that are preventable through primary prevention (such as vaccines against infectious agents or altered lifestyles) may be evaluated from their ERS. For nonhereditary D-tumors, this fraction is high and primary prevention could make a major impact (31). Secondary prevention, obtainable in principle through early detection, could further reduce nonhereditary D-tumor-related deaths and is also instrumental for reducing hereditary D-tumor-related deaths. For R-tumors, primary prevention measures are not likely to be very effective, and secondary prevention should be the major focus.

Here I describe statistical problems with the paper that undermine these conclusions. These problems are so basic that they ought to have been caught in review, but apparently were not.

First, any statistic on "extra risk" should be invariant under a change in the units of time used to measure the rates. Their ERS statistic not only fails to be invariant, but is non-monotonic under a simple change of time units in the data. This renders meaningless all the conclusions based on it.

Second, there is a fundamental misunderstanding of what high correlations imply. The argument that "primary prevention measures are not likely to be very effective" rests on the idea that high correlations between a variable not subject to intervention (number of stem cell divisions) and a target variable (cancer incidence) means that the target variable is mostly non-

susceptible to intervention.

This is wrong on two counts: first, because the correlation is between logarithms, it is possible for a second, unmeasured factor to vary, in this case over four orders of magnitude, and still maintain the correlation of $0.8$ for the data. Second, correlations only put limits on the existing variation of unknown factors; they have nothing to say about novel interventions that may be developed which change that variation. This latter point is well developed in Feldman and Lewontin (1975), an article prompted by the misuse of heritability measures of human intelligence. The potential misuse of correlation measures in decisions about cancer research prompts the work here.

Putting aside the statistical problems in Tomasetti and Vogelstein (2015), further exploration of their data reveal some tantalizing clues to the biology of cancer.

First, the best fit to the data indicates that cancer incidence grows not in proportion to the number of stem cell divisions in a tissue, but in proportion close to the square root. Second, an exploration of different combinations of primary data, $s$ and $d$, they use to estimate the lifetime number of stem cell divisions shows the correlations with the lifetime cancer incidence are maximized almost *precisely* by the formula they use, $\text{LSCD} = s(2 + d) - 2$. This suggests that their estimate for the number of stem cell divisions or a closely related formula is a real biological factor in the incidence of cancer, while not ruling out the possibility of other central factors.

## The Ill-Behaved "Extra Risk Score" (ERS)

Tomasetti and Vogelstein (2015) introduce their ERS statistic:

> We next attempted to distinguish the effects of this stochastic, replicative component from other causative factors—that is, those due to the external environment and inherited mutations. For this purpose, we defined an extra risk score (ERS) as the product of the lifetime risk and the total number of stem cell divisions ($\log_{10}$ values). . . .

> The ERS provides a test of the approach described in this work. If the ERS for a tissue type is high—that is, if there is a high cancer risk of that tissue type relative to its number of stem cell divisions—then one would expect that environmental or inherited factors would play a relatively more important role in that cancers risk (see the supplementary materials for a detailed explanation). It was therefore notable that the tumors with relatively high ERS were those with known links to specific environmental or hereditary risk factors (Fig. 2, blue cluster).

The most straightforward statistic to measure the "cancer risk of that tissue type relative to its number of stem cell divisions" would be the ratio of risk to the number of cell divisions,

LCI/LSCD. On the logarithmic scale this would be $\log_{10} \text{LCI} - \log_{10} \text{LSCD}$. For unknown reasons, the authors instead devise a statistic $\text{ERS}(\text{LCI}, \text{LSCD}) := \log_{10} \text{LCI} \times \log_{10} \text{LSCD}$. They also considered using $\log_{10} \text{LCI} / \log_{10} \text{LSCD}$ (see their Supplement) but reject it, not on first principles, but because it is "suboptimal":

> "Note that using the ratio between the $\log_{10}$ values of $r$ [LCI] and $lscd$, instead of the product, would be sub-optimal to estimate the extra risk. . . . When ERS is defined as the product rather than the ratio, the expected relationship is evident".

One expected relationship that should be evident for a measurement of extra risk is that it be invariant under a change of time units. If the time units were changed from *per lifetime* to *per lifetime*/$T$, this should be irrelevant to a measure of "cancer risk of that tissue type relative to its number of stem cell divisions".

However, we find that ERS is extremely ill-behaved in this regard: it is not invariant, and even worse, it is non-monotonic. This can be seen from its expansion:

$$\begin{aligned} &\text{ERS}(\text{LSCD}/T, \ \text{LCI}/T) \\ &\quad = \text{ERS}(\text{LSCD}, \ \text{LCI}) \\ &\qquad - T(\log_{10} \text{LSCD} + \log_{10} \text{LCI}) + (\log_{10} T)^2. \end{aligned}$$

The relationship between ERS for the data set using time units $T = 1$ and $T = 1000$ is plotted in Figure 1. We see that a simple change of time measurement units essentially scrambles the ERS scores.
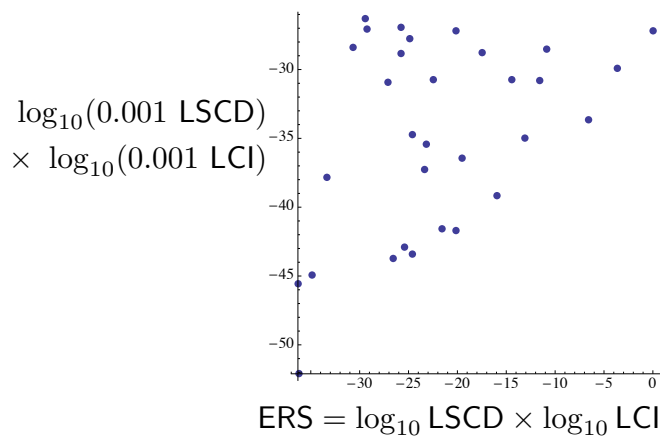


Figure 1: Non-monotonicity of the "Extra Risk Score" (ERS) under change of time units for the cancer incidence data and the number of stem cell divisions. Here the time unit is rescaled by a factor of $1/T = 0.001$.

The non-monotonicity of the ERS statistic in Tomasetti and Vogelstein (2015) under a simple change of time units renders any conclusions based on the ERS meaningless. This includes their Fig. 2, their categorization of "R-tumors" and

"D-tumors", their K-means cluster analysis (Supplement Fig. S2), and therefore the conclusion of their paper (repeated from above):

> The maximum fraction of tumors that are preventable through primary prevention (such as vaccines against infectious agents or altered lifestyles) may be evaluated from their ERS. For nonhereditary D-tumors, this fraction is high and primary prevention could make a major impact ... For R-tumors, primary prevention measures are not likely to be very effective, and secondary prevention should be the major focus.

## Residual Risk

Clearly, Tomasetti and Vogelstein's $\mathsf{ERS} = \log_{10} \mathsf{LCI} \times \log_{10} \mathsf{LSCD}$ is so ill-behaved that it has no meaning. The two natural choices for measuring extra risk are:

1. **Normalized Incidence**: This is simply the lifetime cancer incidence divided by the lifetime number of stem cell divisions. In $\log$ scale:

$$\mathsf{NI} := \log_{10} \mathsf{LCI} - \log_{10} \mathsf{LSCD}.$$

2. **Residual Risk**: This is the lifetime cancer incidence divided by the incidence predicted from the linear regression. In $\log$ scale:

$$\mathsf{RR} := \log_{10} \mathsf{LCI} - \mathsf{P}[\log_{10} \mathsf{LCI}],$$

where

$$\mathsf{P}[\log_{10} \mathsf{LCI}] = 0.533 \log_{10} \mathsf{LSCD} - 7.61 \qquad (1)$$

is the linear predictor function from regression of $\log_{10} \mathsf{LCI}$ on $\log_{10} \mathsf{LSCD}$ (plotted with the data in Figure 2):

In Table 1 the 31 tumor types are sorted by their residual risk values,

$$\mathsf{RR} = \log_{10} \mathsf{LCI} - (0.533 \log_{10} \mathsf{LSCD} - 7.61).$$

(this has been posted on at least one online blog, "Peer 3" (2015)).

Topping the list is lung cancer in smokers, followed principally by the tumors labeled with inherited or viral risk factors. This supports the argument of Tomasetti and Vogelstein (2015) for "incorporation of a replicative component as a third, quantitative determinant of cancer risk" in addition "environmental or inherited factors."

We also see something remarkable. The types identified by Tomasetti and Vogelstein as "excess risk" based on the spurious ERS statistic largely maintain their positions in the top of the list based on the RR statistic. The reason for this is pure coincidence: it happens that, on this data set, the correlation between RR and ERS is fortuitously high: $r(\mathsf{ERS}, \mathsf{RR}) = 0.83$, $\rho(\mathsf{ERS}, \mathsf{RR}) = 0.80$.
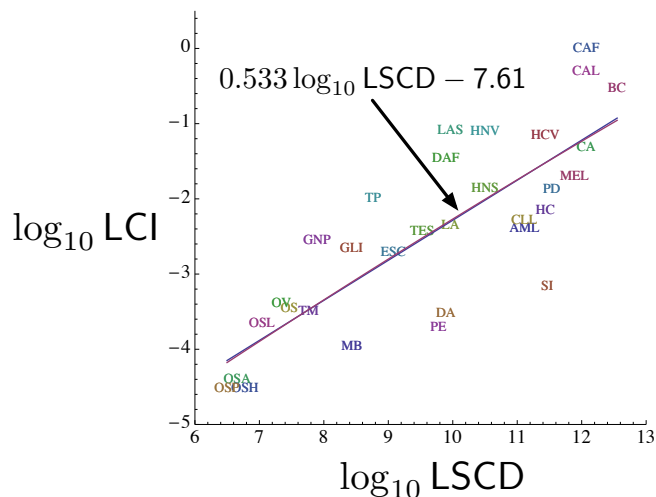


Figure 2: The 31 tissues types (abbreviations) with the linear prediction function.

That this is fortuitous can be seen in the completely different mathematical structures of ERS and RR:

$$\mathsf{ERS} = \log_{10} \mathsf{LCI} \times \log_{10} \mathsf{LSCD},$$
$$\mathsf{RR} = \log_{10} \mathsf{LCI} - (0.533 \log_{10} \mathsf{LSCD} - 7.61).$$

For comparison, the rank correlation between ERS and ERS with time units rescaled by $0.001$ is only $0.24$ (Figure 1), and between ERS and NI it is only $0.09$. Note that $\rho(\mathsf{NI}, \mathsf{RR}) = 0.618$.

This coincidental correlation may have been the reason that ERS was retained by the authors, because it sorted the tumor types in an order similar to the vertical position of the points relative to the predictor line in Figure 2.

If one wanted to distinguish "R-tumors" from "D-tumors" based on Table 1, one certainly could, since there are several large gaps between RR values, notably between Gallbladder and Glioblastoma. But an goodness-of-fit test (Anderson and Darling, 1952) shows the distribution of RR values to be indistinguishable from a Gaussian normal distribution ($P = 0.977$), and the gap sizes not significantly different from an exponential distribution ($P = 0.51$). Caution should therefore be used in making any claims based on these gaps.

We also see in Table 1 that the residual risk spans $2.862 = 1.211 + 1.651$ orders of magnitude. A correlation of $0.8$ between the logarithms of LSCD and LCI therefore does not preclude three orders-of-magnitude variation in LCI due to other factors. This is elaborated upon in the next section. If we consider only those tumor types with negative RR, there are still almost 2 orders of magnitude variation due to unknown factors. If not from errors in the data, something is suppressing the incidence of certain tissue tumors to only 2% of that predicted from the $\log$ LSCD regression.

Table 1: Excess of $\log_{10}$ LCI above the best fit linear regression: RR $= \log_{10}$ LCI $- (0.533 \log_{10}$ LSCD $- 7.61)$. In **bold face** are the types called "D-tumors" in Tomasetti and Vogelstein (2015). Their position near the top of the list is due to the coincidentally high rank correlation between $\log_{10}$ LCI $\times \log_{10}$ LSCD and $\log_{10}$ LCI $- (0.533 \log_{10}$ LSCD $- 7.61)$ on the data set.

| RR | Tumor Type |
|---|---|
| 1.211 | **Lung adenocarcinoma (smokers)** |
| 1.183 | **Colorectal adenocarcinoma with FAP** |
| 0.952 | **Thyroid papillary/follicular carcinoma** |
| 0.916 | **Head & neck squamous cell carcinoma with HPV-16** |
| 0.886 | **Duodenum adenocarcinoma with FAP** |
| 0.882 | **Colorectal adenocarcinoma with Lynch syndrome** |
| 0.853 | Gallbladder non papillary adenocarcinoma |
| 0.460 | Glioblastoma |
| 0.403 | **Basal cell carcinoma** |
| 0.373 | **Hepatocellular carcinoma with HCV** |
| 0.314 | Ovarian germ cell |
| 0.201 | Osteosarcoma of the legs |
| 0.178 | Osteosarcoma |
| 0.156 | Head & neck squamous cell carcinoma |
| 0.106 | Testicular germ cell cancer |
| 0.062 | Esophageal squamous cell carcinoma |
| −0.016 | Thyroid medullary carcinoma |
| −0.045 | Lung adenocarcinoma (nonsmokers) |
| −0.136 | **Colorectal adenocarcinoma** |
| −0.334 | Osteosarcoma of the arms |
| −0.373 | Osteosarcoma of the pelvis |
| −0.400 | Pancreatic ductal adenocarcinoma |
| −0.411 | Melanoma |
| −0.520 | Osteosarcoma of the head |
| −0.593 | Chronic lymphocytic leukemia |
| −0.627 | Hepatocellular carcinoma |
| −0.696 | Acute myeloid leukemia |
| −0.840 | Medulloblastoma |
| −1.181 | Duodenum adenocarcinoma |
| −1.312 | Pancreatic endocrine (islet cell) carcinoma |
| −1.651 | Small intestine adenocarcinoma |
| 0.000 | **Total** |

# The fallacy that *high correlations preclude intervention*

Tomasetti and Vogelstein point out that the correlation between $\log_{10}$ LSCD and $\log_{10}$ LCI is extremely robust to noise, and therefore makes their results robust. What they do not realize is that this very robustness argues against their conclusion that "for R-tumors, primary prevention measures are not likely to be very effective", because it allows large variation in unknown factors that could also control cancer incidence rates while having little effect on the correlation.

Tomasetti and Vogelstein measure this robustness by adding noise to the data to see how much it changes the correlation. They add both Gaussian and uniformly distributed random noise to their estimates of lifetime number of stem cell divisions. For the uniform variation, they examine the correla-

tions between LCI and LSCD $+ \mathcal{U}(-2, 2)$, where $\mathcal{U}(-2, 2)$ is a uniformly distributed random variable on the interval $[-2, 2]$. Under 10,000 replicates they find the addition of this four orders-of-magnitude noise only reduces the median value of Spearman's rank correlation $\rho$ (Spearman, 1904) from $0.81 = \rho(\text{LCI}, \text{LSCD})$ to $0.67 = \rho(\text{LCI}, \text{LSCD} + \mathcal{U}(-2, 2))$, which is still significantly different from zero:

> Thus, though the total range for LSCD is $\sim 6$ orders of magnitude and we allowed four 4 [sic] orders of magnitude variation for each data point, the correlations generated were always statistically significant. This provides strong evidence that our results are robust. [Supplement p. 11]

But the robustness of the high correlation to order-of-magnitude variations in the data has its converse implication: it means that high correlation cannot rule out order-of-magnitude variation in other factors that may determine the rates of cancer. We see this in Table 1.

To illustrate this converse implication, suppose that LCI were determined by two factors, LSCD, and another, unknown preventative factor "X" that reduces the rate of cancer in proportion to its value. How much variation in X could there be and still obtain correlations of $0.8$ between $\log_{10}$ LSCD and $\log_{10}$ LCI?

The general problem is to compute the correlation, $r$ (Pearson, 1901), between $Y$, and $Y$ plus uncorrelated noise, $Z$:

$$r(Y, Y + Z) = \frac{\text{Cov}(Y, Y + Z)}{\sqrt{\text{Var}(Y)}\sqrt{\text{Var}(Y + Z)}}$$
$$= \frac{1}{\sqrt{1 + \frac{\text{Var}(Z)}{\text{Var}(Y)}}}.$$

The solution of $r(Y, Y + Z) = 0.8$ is $\text{Var}(Z) = 0.5625\,\text{Var}(Y)$.

To be concrete, let LCI $= c * \text{LSCD} * \text{X}$, where the constant $c$ does not enter into the correlation. Then $\log_{10}$ LCI $= \log_{10} c + \log_{10}$ LSCD $+ \log_{10}$ X. We let $\log_{10}$ LSCD be distributed as a uniform random variable on the interval $[6.50, 12.55]$, the actual range from the data. We let $\log_{10}$ X be distributed as an independent random variable uniform on $[-w, 0]$, and ask how big $w$ can be and yet maintain $0.8 = r(\log_{10} \text{LCI}, \log_{10} \text{LSCD})$. The solution is $w = 4.5$—four orders of magnitude.

It may be helpful to see what the distribution of this hypothetical preventative factor X looks like, in Figure 3. We see huge variation in the population incidence of cancer due to hypothetical factor X, even though there is still a correlation of $0.8$ between $\log_{10}$ LSCD and $\log_{10}$ LCI. Therefore, the correlation of $0.8$ does not rule out the presence of other factors—genetic, environmental, physiological, or measurement error—that can produce over four orders-of-magnitude variation in cancer rates. Such factors could potentially include environmental variation that is subject to intervention.
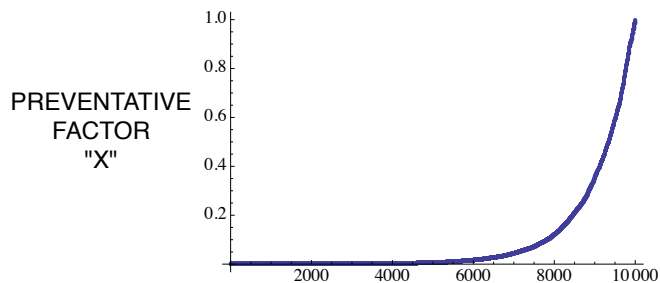
Figure 3: Variation in a hypothetical multiplicative "preventative factor X" on cancer incidence consistent with a correlation of 0.8 between $\log_{10}$ LSCD and $\log_{10}$ LCI. Sorted along the X-axis are 10,000 samples of X, where $\log_{10}$ X $\sim U[-4.5, 0]$.

## Further Explorations

The problems described above do not diminish the fact that Tomasetti and Vogelstein (2015) have produced an important data set, and the strong correlation they find between the logarithms of lifetime number of stem cell divisions and the rates of cancer in different tissues provides evidence for a real biological phenomenon. Their data therefore deserve deeper investigation.

### A One-Half Power Law?

The predictor function (2) translates to an approximate power law of

$$\text{LCI} \approx 10^{-7.61} \text{ LSCD}^{0.533}. \tag{2}$$

A power law with exponent 0.533 is not what we would have expected from a simple stem cell division hypothesis. The simple hypothesis would predict the cancer incidence should grow *in proportion* to the number of stem cell divisions, with exponent 1.00. What we find instead is that it grows in proportion to approximately the *square root* of the number of stem cell divisions (this has been noted in comments on at least one online blog, Kuenzel (2015)).

The value 0.533 may, however, reflect *regression dilution*—a systematic reduction in the slope of the predictor function due to large amounts of noise in the independent variable. This is likely here because the estimate of the number of stem cell divisions, LSCD, combines multiple and sometimes uncertain primary data. Improvements in the estimation of LSCD and inclusion of more tissue types in the data would therefore be expected to increase the slope.

If, however, further data lent support to a power law with an exponent of one-half, it would invite biological speculation. One would want to examine probabilistic models for the mutational basis of carcinogenesis to see where probabilities scaled in proportion to the square root of the number of cell divisions. Random walks and diffusions provide another possibility, since

the mean squared displacement grows in proportion to time or the number of steps taken. In such cases we would need to ask what properties of cells follow a random walk under cell division. Models of telomere length come to mind here (Blythe and MacPhee, 2013; Duc and Holcman, 2013). Another source could be geometry—boundaries of regions that grow along two-dimensions (e.g. the perimeter of a growing circle, or the surface of an elongated cylinder growing in diameter) can have close to a square-root relationship to that growth. Provided the right geometry, a square-root relationship could conceivably emerge if the boundaries of tissues played a role in carcinogenesis. There may be plausible sources of one-half power laws to be elicited from the work on scaling relationship in biology (c.f. Savage et al. (2013)).

### The Formula for LSCD

The appearance of the exponent 0.533 in the relationship between LCI and LSCD prompts us to examine more closely how the estimate of LSCD was made. Their Supplement explains that the formula for LSCD is:

$$\text{LSCD} = s(2 + d) - 2, \tag{3}$$

where

$s$ is the total number of stem cells found in a fully developed tissue, and

$d$ is the number of further divisions of each stem cell in the lifetime of that tissue once the tissue is fully developed, due to normal tissue turnover.

Suppose instead of this formula, we explored the space of possible combinations of the primary data $s$ and $d$ (i.e we are pursuing *symbolic regression* (Koza, 1990)). Could modern data mining software discover the stem cell division hypothesis by finding patterns in the data (c.f. Schmidt and Lipson (2009))? Could we find combinations of $s$ and $d$ with higher correlations to LCI?

We first explore the effect of changing the additive constants in (3). For Spearman's $\rho$, the constant $-2$ is irrelevant. We embed $s(2+d)-2$ in the family of formulae, $\psi(c) := s(c+d) - 2$. For $c \in [0, 20]$, $\rho(\psi(c), \text{LCI})$ varies only between 0.776 and 0.810, with a broad peak for $0.8 \le c \le 6.1$. The robustness of $\rho$ to variation in the constants means the data give no particular validation to their values in the formula for LSCD (3).

Since $\log_{10}$ LSCD $\approx \log_{10} s + \log_{10} d$, hence (1) is approximated by

$$\text{P}[\log_{10} \text{LCI}] \approx 0.533 \left(\log_{10} s + \log_{10} d\right) - 7.61.$$

Suppose instead of equal weights on $\log_{10} s$ and $\log_{10} d$, we embed $s(2+d) - 2$ in the family of formulae,

$$\phi(t) := s^{(2-t)}(2 + d^t) - 2.$$

The parameter $t$ varies $\phi(t)$ continuously from $\phi(0) = 3s^2 - 2$, to $\phi(1) = \mathsf{LSCD}$, to $\phi(2) = d^2$. We find that

$$\rho(\phi(0), \mathsf{LCI}) = \rho(s, \mathsf{LCI}) = 0.68,$$
$$\rho(\phi(1), \mathsf{LCI}) = \rho(\mathsf{LSCD}, \mathsf{LCI}) = 0.81,$$
$$\rho(\phi(2), \mathsf{LCI}) = \rho(d, \mathsf{LCI}) = 0.60.$$



Figure 4: Correlations $r(\log_{10}\mathsf{LCI}, \log_{10}\phi(t))$ and $\rho(\mathsf{LCI}, \phi(t))$ for $\phi(t) := s^{(2-t)}(2 + d^{\,t}) - 2$. The peaks are near $t = 1$, where $\phi(1) = \mathsf{LSCD}$.

For the range of $t$, Figure 4 plots the correlations $\rho(\phi(t), \mathsf{LCI})$ and $r(\phi(t), \mathsf{LCI})$. There peaks are very near the value $t = 1$ where $\phi(1) = \mathsf{LSCD}$.

Therefore, by simply looking for correlations between $\mathsf{LCI}$ and different combinations of $s$ and $d$, we come up with a formula very close to the formula developed by Tomasetti and Vogelstein *from first principles* for the total number of stem cell divisions in a tissue. There is no *a priori* reason that the peak should occur near $t = 1.0$, since $s$, $d$, and $\mathsf{LCI}$ are independently obtained data sets (the rank correlation between $s$ and $d$ is only 0.23).

Could this be a coincidence? We can get some idea of that likelihood. Bootstrap resampling of the 31 tissue types for the values of $t$ that maximize $\rho(\phi(t), \mathsf{LCI})$ gives a central 50% interval of $(0.909, 1.009)$, with a median and sharp mode at $t = 0.971$. Subjecting $s$ and $d$ to multiplicative noise (replacing $s$ and $d$ by $s \times \nu_1$ and $d \times \nu_2$, where $\ln \nu_1, \ln \nu_2 \sim \mathcal{N}(0, 1/8^2)$, $\nu_1, \nu_2$ independent), gives a median of 0.954 and 50% central range of $(0.928, 0.973)$.

So, yes it could be a coincidence that the optimal $t$ is so close to 1.0. Chance can't be ruled out as to why the formula of Tomasetti and Vogelstein seems to provide the maximal correlation to the rates of cancer incidence. We cannot rule out there being other relationships between $s$ and $d$ that give even higher correlations with the cancer incidence data. It will require expansion of the data set to more tissue types, more precise estimates of the number of stem cell divisions, and exploration of other models of the biology to resolve this question.

Nonetheless, the optimality of the formula of Tomasetti and Vogelstein is at least suggestive: that something close to $s \times d$ has a real biological role in the incidence of cancer. It is a separate form of evidence from the 0.8 magnitude of the correlation itself.

## Discussion

The publication of Tomasetti and Vogelstein (2015) was announced by Johns Hopkins University with the press release headline, "Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows." Probability theory was created to make reasoning about "luck" a rigorous science. People's reasoning about probabilities—luck—is a notoriously error-prone activity, and the characterization of these errors is now a scientific field that has burgeoned since the pioneering work of (Tversky and Kahneman, 1974). There is a wide intuition that there is a zero-sum tradeoff between control and luck in determining events: the more control we have, the less we depend on luck. This intuition of a tradeoff resonates with the mathematical structure of the analysis of variance, in which total variance is partitioned into a sum of variances and covariances. However, when two factors, "luck" ($L$) and "control" ($C$) interact multiplicatively, as $L \times C$, then there is no necessary tradeoff between them at all. And as described here, because of the nature of logarithms, a high correlation between $\log L$, and $\log L + \log C$, does not preclude large variation in $C$.

Beyond this issue of correlations between logarithms, and the erroneous use of the ill-behaved "extra risk score" ($\mathsf{ERS}$), there is the additional pitfall of cognitive framing. If Tomasetti and Vogelstein had presented an 80% correlation between logarithms of the number of stem cell divisions and cancer *mortality* for different tissues, instead of *incidence*, would they have then argued "cancer treatment measures are not likely to be very effective, and prevention should be the major focus"? Of course not, because novel cancer treatments are something new under the sun. The cognitive frame in this case draws on our familiarity with the history of medicine in which many new cures have been found that obliterate past correlations.

The experience of "dramatic cure" has no parallel in experiences of "dramatic prevention" (except to a statistician) because prevention is a non-event. Therefore, the fallacy that "high correlation precludes intervention" can find an easier home when reasoning about prevention.

Cures, primary prevention, and secondary prevention are simply interventions at different stages of disease. Correlations found in the present do not bear on what interventions may be found in the future for any of these disease stages. If anything, the track record of human discovery to date has been that cancer prevention is much easier to discover than cancer cure. But the future has yet to be written.

Leaving aside the erroneous statistics of Tomasetti and Vogelstein (2015) and the conclusions based on them, as *biology* it is a significant finding that there is a high rank correlation

between (1) cancer incidence, and (2) the estimated number of stem cell divisions in a tissue. Additional upport for that significance is presented here: that among various possible combinations of the primary data, the one that produces the highest correlation to cancer incidence is precisely the formula Tomasetti and Vogelstein obtain from biological first principles. Although this could be a statistical coincidence and requires additional data to confirm, it can be seen as a novel form of support for the hypothesis stem cell divisions are an independent causal factor for cancer.

**Materials and Methods**

Data sets were obtained from the Supplement to Tomasetti and Vogelstein (2015). All computations and statistics were carried out using *Mathematica*™. The code used is available upon request.

# Literature Cited

Anderson, T. W. and Darling, D. A. 1952. Asymptotic theory of certain" goodness of fit" criteria based on stochastic processes. *The Annals of Mathematical Statistics*, pages 193–212.

Blythe, R. A. and MacPhee, C. E. 2013. The life and death of cells. *Physics*, 6:129.

Duc, K. D. and Holcman, D. 2013. Computing the length of the shortest telomere in the nucleus. *Physical Review Letters*, 111(22):228104.

Feldman, M. W. and Lewontin, R. C. 1975. The heritability hang–up. *Science*, 190(4220):1163–1168.

Koza, J. R. 1990. Genetic programming: A paradigm for genetically breeding populations of computer programs to solve problems. Technical Report CS-TR-90-1314, Stanford University, Department of Computer Science, Stanford, CA.

Kuenzel, M. 2015. Thoughts on "are two thirds of cancers really due to bad luck?". *The Stats Guy (blog)*, January 3 (3:58 pm). `http://www.statsguy.co.uk/are-two-thirds-of-cancers-really-due-to-bad-luck/`.

Pearson, K. 1901. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572.

"Peer 3" 2015. Comments. *Pub Peer (blog)*, January 3 (3:54 pm UTC). `http://pubpeer.com/publications/867D3AFB2EBDF34A37FFCBA075D1BE`.

Savage, V. M., Herman, A. B., West, G. B., and Leu, K. 2013. Using fractal geometry and universal growth curves as diagnostics for comparing tumor vasculature and metabolic rate with healthy tissue and for predicting responses to drug therapies. *Discrete and Continuous Dynamical Systems. Series B*, 18(4).

Schmidt, M. and Lipson, H. 2009. Distilling free-form natural laws from experimental data. *Science*, 324(5923):81–85.

Spearman, C. 1904. The proof and measurement of association between two things. *The American Journal of Psychology*, 15(1):72–101.

Tomasetti, C. and Vogelstein, B. 2015. Variation in cancer risk among tissues can be explained by the number of stem cell divisions. *Science*, 347(6217):78–81.

Tversky, A. and Kahneman, D. 1974. Judgment under uncertainty: Heuristics and biases. *Science*, 185(4157):1124–1131.